

# Binary Compressed Imaging

Aurélien Bourquard, *Student Member, IEEE*, and Michael Unser, *Fellow, IEEE*

**Abstract**—Compressed sensing can substantially reduce the number of samples required for conventional signal acquisition, at the expense of an additional reconstruction procedure. It also provides robust reconstruction when using quantized measurements, including in the one-bit setting. In this paper, our goal is to design a framework for binary compressed sensing that is adapted to images. Accordingly, we propose an acquisition and reconstruction approach that complies with the high dimensionality of image data and that provides reconstructions of satisfactory visual quality. Our forward model describes data acquisition and follows physical principles. It entails a series of random convolutions performed optically followed by sampling and binary thresholding. The binary samples that are obtained can be either measured or ignored according to predefined functions. Based on these measurements, we then express our reconstruction problem as the minimization of a compound convex cost that enforces the consistency of the solution with the available binary data under total-variation regularization. Finally, we derive an efficient reconstruction algorithm relying on convex-optimization principles. We conduct several experiments on standard images and demonstrate the practical interest of our approach.

**Index Terms**—Acquisition devices, bound optimization, compressed sensing, conjugate gradient, convex optimization, inverse problems, iteratively reweighted least squares, Nesterov’s method, point-spread function, preconditioning, quantization.

## I. INTRODUCTION

In the context of compressed sensing, the amount of data to be acquired can be substantially reduced as compared to conventional sampling strategies [1]–[6]. The key principle of this approach is to compress the information before it is captured, which is especially beneficial when the acquisition process is expensive in terms of time or hardware. For instance, in their previous work [7], Boufounos *et al.* investigated the performance of compressed sensing in the binary case where the extreme coarseness of the quantization must typically be compensated by taking more numerous measurements than in the classical case. The original signal can then be recovered from the available measurements through numerical reconstruction, whose computational complexity exhibits a strong dependence on the structure of the forward model. Consequently, specialized acquisition approaches are required for compressed sensing when dealing with large-scale data such as images. For instance, we were able to extend in [8] the central principles of [7] to image acquisition and reconstruction. Our associated forward model generates binary measurements that are based on random-convolution principles [1]. Though demonstrating satisfactory reconstruction capability for image data, this method tends to create

spatial redundancy in the associated measurements, which is suboptimal from the perspective of information content.

In this paper, our first contribution is to propose a general framework for the binary compressed sensing of images. Based on [8], we devise an extended forward model that can take several binary captures of a given grayscale image. Each of these acquisitions corresponds to one distinct convolution performed by an optical system. The flexibility of our approach allows us to improve the statistical properties of the associated binary data, which ultimately increases the quality of reconstructions.

Using a variational formulation to express our reconstruction problem, our second contribution is a fast reconstruction algorithm that uses bound-optimization principles. The design of this algorithm bears similarities with the optimization techniques [9], [10]. It yields an iteratively reweighted least-squares (IRLS) procedure that is easily parameterized and that converges in few iterations.

We introduce our forward model for image acquisition in Section II. In Section III, we express our reconstruction problem as the minimization of a compound cost functional. Based on convex optimization, we derive the reconstruction algorithm in Section IV. In Section V, we perform several experiments on standard grayscale images. We extensively discuss them and conclude our work in Section VI.

## II. FORWARD MODEL

### A. General Structure

In this section, we establish a convolutive physical model that generates  $L$  binary measurement sequences  $\gamma_i$  from a given two-dimensional (2D) continuously defined image  $f$  of unit square size. Following a design similar to the one of [8], each of these sequences is obtained through optical convolution of  $f$  with a distinct pseudo-random filter  $h_i$  followed by acquisition through binary sensors.

Specifically, each convolved image  $f * h_i$  is sampled and binarized by a uniform 2D CCD-like array of  $M_0 \times M_0$  sensors, the specific form of  $h_i$  being defined in Section II-B. The actual sampling process is regular but nonideal, meaning that each sensing area of side  $M_0^{-1}$  has some pre-integration effect modeled by some spatial filter  $\phi$ . Therefore, the global convolutive effect of our model before sampling corresponds to the spatial kernels  $\chi_i^0 = h_i * \phi$ , yielding the pre-filtered intermediate images

$$f_i^0(\mathbf{x}) = (f * \chi_i^0)(\mathbf{x}), \quad (1)$$

where the vector  $\mathbf{x} \in \mathbb{R}^2$  denotes the 2D spatial coordinates. Then, the sensor array samples each image  $f_i^0$  with a step  $T = M_0^{-1}$ , which produces the sequences  $g_i^0$  defined for each index  $\mathbf{k} \in \mathbb{Z}^2$  as

The authors are with the Biomedical Imaging Group (BIG), École polytechnique fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland (email: aurelien.bourquard@epfl.ch; michael.unser@epfl.ch).

Digital Object Identifier xx.xxxx/TIP.2011.xxxxx

$$g_i^0[\mathbf{k}] = f_i^0(\mathbf{x})|_{\mathbf{x}=\mathbf{k}T}. \quad (2)$$

Unlike [8], we allow for a finite-differentiation process to take place before the final quantization step. Denoting the corresponding discrete filters as  $\nabla_i$ , the non-quantized measurements  $g_i$  are obtained as

$$g_i[\mathbf{k}] = (g_i^0 \star \nabla_i)[\mathbf{k}], \quad (3)$$

where  $\star$  denotes a discrete convolution. These operations can be efficiently performed by the sensor array itself, for instance using voltage comparators. As discussed in the experimental section, finite differentiation brings improvements in terms of reconstruction quality and simplifies the calibration of the system. Note that no finite differentiation occurs when taking  $\nabla_i$  to be<sup>1</sup> the discrete unit sample  $\delta[\cdot]$ .

Defining  $\tau$  as a common threshold value, the quantized measurement sequences  $\gamma_i$  are obtained as

$$\gamma_i[\mathbf{k}] = \begin{cases} +1, & g_i[\mathbf{k}] \geq \tau \\ -1, & \text{otherwise.} \end{cases} \quad (4)$$

The measurements  $\gamma_i$  can be selectively stored according to discrete spatial indicator functions  $\omega_i$ . Each  $\gamma_i[\mathbf{k}]$  is actually kept and counted as a measurement if and only if the value  $\omega_i[\mathbf{k}] \in \{0, 1\}$  is unity for the same  $\mathbf{k}$ . Note that, before binarization, every measurement  $g_i$  is a mere linear functional of  $f$ .

The successive operations that are involved in our forward model simplify to one single convolution in the continuous domain without subsequent discrete filtering, as summarized in Figure 1. The equivalent spatial impulse response  $\chi_i$  of the filter corresponds to  $\chi_i(\mathbf{x}) = \sum_{\mathbf{k}} \nabla_i[\mathbf{k}] \chi_i^0(\mathbf{x} - T\mathbf{k})$ . To sum up, our forward model yields  $M = \Lambda LM_0^2$  binary measurements of the continuously defined image  $f$  in the form of  $L$  distinct binary sequences, where  $\Lambda$  is the storage ratio associated to the functions  $\omega_i$ . These captured sequences are complementary, as they are associated with distinct random convolutions before sampling, binarization, and masking through the  $\omega_i$ . Since the latter process allows to decrease  $M$ , the resolution  $M_0$  can be kept constant. This avoids high-frequency losses due to coarse-sensor integration.

Besides reducing data storage, the process of binary quantization potentially consumes far less power than standard analog-to-digital converters, and is less susceptible to the nonlinear distortion of analog electronics [9]. Binary sensors are also associated with very high sampling rates in general [7]. In that regard, the selective subsampling that we specify by  $\omega_i$  may also lead to further reductions of the acquisition time if fewer measurements are required; the acquisition of the selected samples can indeed be performed efficiently through randomly addressable image sensors<sup>2</sup>.

<sup>1</sup>The symbol  $\cdot$  denotes a dummy variable. It can be used to create new function definitions based on existing ones. For instance,  $f(\cdot - \mathbf{k})$  corresponds to the original image  $f$  shifted spatially by  $\mathbf{k}$ .

<sup>2</sup>Image sensors that are based on the complementary-metal-oxide-semiconductor (CMOS) technology allow for parallel and random access, as opposed to other architectures that can only perform sequential readout [11]. While the potential benefits of binary sensors further motivate our work, the proper development of such elements for optics remains to be addressed.

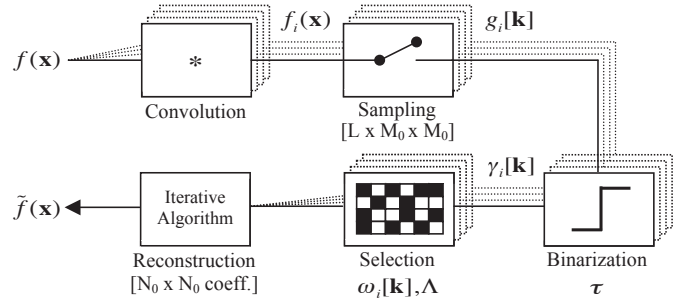


Fig. 1. General framework. The unknown continuously defined image  $f$  is first convolved with  $L$  distinct kernels  $\chi_i$ , producing the intermediate images  $f_i = f \star \chi_i$ . Each  $f_i$  is then sampled with step  $T$  to obtain the sequences  $g_i$ . The last acquisition step consists in pointwise binarization with threshold  $\tau$ , resulting in the binary measurements  $\gamma_i$ . When retained, the latter constitute the available information on the original data. Based on these selected measurements, and assuming that the forward model is known, our reconstruction algorithm produces an estimate  $\tilde{f}$  of the original image.

### B. Pseudo-Random Optical Filters

As mentioned in Section II-A, the  $L$  filters  $h_i$  are associated to optical convolution operations. Accordingly, we make each  $h_i$  correspond to a distinct spatially invariant point-spread function (PSF) that is generated by the same optical model. In our setup shown in Figure 2, the image  $f$  is associated with light intensities defined on a plane. For each of the  $L$  acquisitions, the intensities measured by the sensor array after optical propagation correspond to the convolution  $f \star h_i$  up to geometrical inversion.

The specific form of  $h_i$  depends on the profile of the central plane of the system called the *Fourier plane* [12]. In our model, this plane transmits light through a circular area and is further equipped for each acquisition with one distinct instance of a phase-shifting plate whose effect is to multiply the transmitted-light amplitudes with pseudorandom phase values. The resulting profile  $q_i$  is modeled as a complex-valued function expressed in normalized spatial coordinates [12].

Considering phase functions  $\mu_i$  composed of square zones, each zone associated with either a  $0$  or  $\pi$  phase shift, we obtain

$$\mu_i(\boldsymbol{\xi}) = \sum_{\mathbf{k} \in \mathbb{Z}^2} \nu_i[\mathbf{k}] \text{rect}(\boldsymbol{\xi} - \mathbf{k}), \quad (5)$$

where the phase values  $\nu_i$  are independent and uniformly distributed random variables from the pair  $\{0, \pi\}$ , where  $\text{rect}$  is the 2D rectangle function, and where  $\boldsymbol{\xi}$  denotes normalized spatial coordinates. The phase-shifting plates associated with the  $\mu_i$  are of finite extent since they only operate inside the transmissive circular area of Figure 2. The latter is designed such that the diameter of the circle covers  $K$  phase zones in the horizontal or vertical direction. It is thus specified by the function

$$\text{circ}(\boldsymbol{\xi}) = \begin{cases} 1, & \|\boldsymbol{\xi}\| \leq K/2 \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

The profile  $q_i$  combines the phase shifts of (5) with the transmissivities of (6). It is defined as

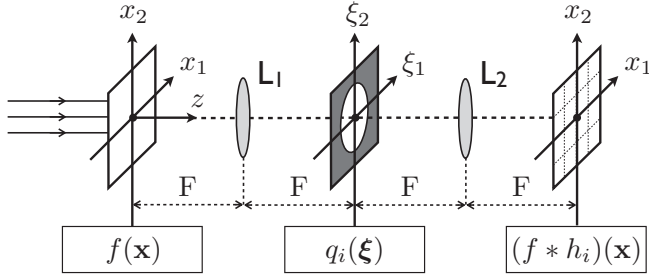


Fig. 2. Optical setup. In our optical model, the image  $f$  maps to light-intensity values. Our optical device transforms this initial image wavefront using elements that are spaced by the same distance  $F$ . Following the direction  $z$  of light propagation, this system called 4F consists in the left plane where the image  $f$  lies, one first lens  $L_1$  of focal length  $F$ , the central plane, one second lens  $L_2$  identical to  $L_1$ , and the last plane containing the propagated wavefront to be captured by the sensor array.

$$q_i(\xi) = \text{circ}(\xi) \exp(-j\mu_i(\xi)). \quad (7)$$

Due to the 4F placement of the lenses, the propagation of light implements a continuous Fourier transform [12]. The light amplitudes are also modulated by  $q_i$  in the Fourier plane. Accordingly, the impulse response of the system is defined up to scale as

$$h_i(\mathbf{x}) = |\mathcal{F}\{q_i\}(\mathbf{x})|^2, \quad (8)$$

where  $\mathcal{F}$  denotes the Fourier transform  $\mathcal{F}\{q_i\}(\mathbf{x}) = \int_{\mathbb{R}^2} q_i(\xi) \exp(-j\mathbf{x}^T \xi) d\xi$ . The use of spatially incoherent illumination<sup>3</sup> and the fact that the measured quantities are light intensities results in a squared modulus in (8). Each filter  $h_i$  is thus nonnegative, and depends upon the corresponding  $\mu_i$  defined in (5). The latter can be generated electronically by a spatial light modulator [1].

### C. Connection with Compressed Sensing

The use of phase masks in our forward model produces random-like patterns in each of our binary-measurement sequences. This closely relates our method to the compressed-sensing paradigm of [7] and requires us to express all our unknowns in discrete form. To this end, we model the continuously defined function  $f$  as the expansion

$$f(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^2} c[\mathbf{k}] \beta^m(\mathbf{x} - \mathbf{k}), \quad (9)$$

where the sequence  $c$  corresponds to  $(N_0 \times N_0)$  real coefficients placed on a regular grid, and where  $\beta^m(\mathbf{x}) = \beta^m(x_1) \beta^m(x_2)$  for  $\mathbf{x} \in \mathbb{R}^2$  is the separable 2D B-spline of degree  $m$ . Given their small support and polynomial-reproduction properties, B-splines are especially adapted from both approximation and computational viewpoints. They thus constitute a suitable approach to represent continuous images [13].

<sup>3</sup>Spatial incoherence means that the phases of the initial wavefront on the left plane of Figure 2 vary with time in uncorrelated fashions. This implies that the effective response of our optical system is linear in intensity rather than in amplitude [12].

Moreover, since our continuous image is modeled as a linear combination of B-spline basis functions, it is equivalently described through the corresponding coefficients. Then, the substitution of (9) into our physical forward model naturally leads to a linear and discrete dependency between the image coefficients and the measurements before quantization. Accordingly, the general relation between the unknown sequence  $c$  and the sequences  $g_i$  can be summarized into the *measurement matrix*  $\mathbf{A} \in \mathbb{R}^{M \times N}$ , whose structure is induced from our continuous-domain formulation. In this paper, vectors refer to lexicographically ordered versions of the corresponding sequences. Using this convention, we obtain

$$\mathbf{g} = \mathbf{A} \mathbf{c}, \quad (10)$$

where  $\mathbf{c}$  contains  $N = N_0^2$  coefficients and where  $\mathbf{g}$  contains  $M$  measurements. Our measurement matrix generalizes [8] and vertically concatenates several terms  $\mathbf{A}_i$  of similar structure as  $\mathbf{A} = (\mathbf{A}_1, \dots, \mathbf{A}_i, \dots, \mathbf{A}_L)$ . These terms are associated with the sequences  $g_i$ . They depend from the corresponding kernels  $\chi_i$  and from the rational sampling step  $T$ . They are defined as

$$\mathbf{A}_i = \Omega_i \mathbf{D}_N \mathbf{B}_i \mathbf{U}_M, \quad (11)$$

where  $\mathbf{D}_i$  and  $\mathbf{U}_j$  denote downsampling-by- $i$  and upsampling-by- $j$  matrices. The integers  $\mathcal{M}$  and  $\mathcal{N}$  are such that the right-hand side of the equality  $M_0/N_0 = \mathcal{M}/\mathcal{N}$  is in reduced form. Given periodic boundary conditions, the circulant matrix  $\mathbf{B}_i$  is associated with the discrete impulse response

$$b_i[\mathbf{k}] = \frac{N}{\mathcal{N} \mathcal{M}_0^2} \left( \chi_i \left( \frac{N \cdot}{\mathcal{N} \mathcal{M}_0^2} \right) * \beta^m \left( \frac{\cdot}{\mathcal{M}} \right) \right) (\mathbf{x})|_{\mathbf{x}=\mathbf{k}}. \quad (12)$$

Finally, each matrix  $\Omega_i$  is linked to  $\omega_i$ . Specifically, it corresponds to an identity matrix whose rows associated with the discarded measurements are suppressed, if any. The overall structure of  $\mathbf{A}$  will prove to be beneficial for the reconstruction in terms of computational complexity. The measurements are indeed related to the coefficients by mere discrete Fourier-transform and resampling operations.

When the unknown vector  $\mathbf{c}$  is sufficiently *sparse* in some adequate basis, which does not need to be known explicitly, the theory of compressed sensing offers guarantees on the quality of reconstruction in terms of robustness to measurement loss or quantization [6], [7], provided that the measurement matrix is appropriate. In the general case, a common and suitable criterion for  $\mathbf{A}$  is to be *statistically incoherent* with any fixed signal representation, which means that the bases of the measurement and sparse-representation domains of the signal are uncorrelated with overwhelming probability [2]. This property has been shown theoretically to strictly hold for matrices consisting of independent and identically distributed (iid) Gaussian random entries [3], [6], and also to nearly hold for other random-matrix ensembles [1], [4], [5], [14], [15]. In this work, we resort to an experimental validation of our measurement matrix for binary compressed sensing. In particular, we shall demonstrate in Section V that our model is suitable for the reconstruction of images from few data, and



that the quality of the solution is linked to relatively simple criteria relying on the measurements themselves.

The appropriateness of  $\mathbf{A}$  in our generalized model is tied to the set of discrete filters  $b_i$  defined in (12) and associated with the matrix terms  $\mathbf{A}_i$ . Indeed, they share similarities with the Romberg's random-convolution pulses proposed in [1] for compressed sensing. Firstly, their discrete Fourier coefficients also have phase values that are randomly distributed in  $[0, 2\pi)$ , given their relation with the profiles (5). Secondly, despite not being strictly all-pass as in [1], our filters are also spread-out in the spatial domain. Due to these properties, the form of  $b_i$  has been shown to yield satisfactory reconstructions in the binary case [8]. Besides being adequate individually, these filters also produce  $L$  distinct sequences  $g_i$  from the same image  $f$  because they are associated with  $L$  distinct pseudorandom phase-mask profiles in (7). In some sense, our multi-acquisition framework is the reverse of multichannel compressed-sensing architectures where one single output sequence combines several source signals through distinct modulation or filtering operations [16], [17]. As will be discussed in Section V, the subsequent thresholding operation (4) that is applied in our method yields binary measurements that follow an equiprobable distribution, as in [8]. The proper specification of the additional acquisition parameters of our system (including  $\omega_i$  and  $L$ ) will allow us to maximize the reconstruction performance while maintaining a high computational efficiency.

### III. FORMULATION OF THE RECONSTRUCTION PROBLEM

For the general problem of binary compressed sensing, the authors of [9] have recently proposed a reconstruction technique that is based on binary iterative hard thresholding (BIHT), using the non-convex constraint that the solution signal lies on the unit sphere. This approach extends previous works [7], [18], and achieves better performance. The work of [10] uses a distinct strategy by formulating a convex reconstruction problem solvable by linear programming. An extension of this principle to the case of noisy measurements is also considered by the same authors in [19].

In this paper, we propose to formulate our image-reconstruction problem in a variational framework. Specifically, our solution is expressed as the minimum of a convex functional that includes data-fidelity and regularity constraints. Using bound-optimization principles, the convexity of this functional is exploited in Section IV to derive an efficient iterative-reconstruction algorithm. The latter can handle large-scale problems because, from a computational perspective, it involves the application of the forward model (whose form is essentially convolutive in our case) and of its adjoint inside each iteration as in other methods. Furthermore, besides quality considerations, the specific structure of our reconstruction problem will allow us to maximize iterative performance through preconditioning and Nesterov's acceleration [20].

The available data consist of the measurements  $\gamma_i$  obtained according to Section II. In addition, we suppose that  $\mathbf{A}$  is known. Its components can be deduced physically from the  $L$  impulse responses  $h_i$  produced by the optical system, or,

more indirectly, from the phase-mask profiles  $\mu_i$ . Based on that information, our goal is to reconstruct an accurate continuously defined estimate  $\tilde{f}$  of the original image  $f$  according to some sparsity prior  $\mathcal{R}$ . Specifically, we demand our reconstructed coefficients  $\tilde{c}$  to minimize

$$\mathcal{J}(\tilde{c}) = \mathcal{D}(\tilde{c}) + \lambda \mathcal{R}(\tilde{c}). \quad (13)$$

The first scalar term  $\mathcal{D}$  imposes the fidelity of the solution to the known binary measurements  $\gamma_i$ . Due to quantization, fidelity alone is in general under-constrained and accurate only up to contrast and offset. Then, the regularization term  $\mathcal{R}$ , weighted by  $\lambda$ , encourages the sparsity of the reconstruction.

#### A. Data Term

The role of our data-fidelity constraint is to ensure that the reintroduction of the reconstructed continuously defined image  $\tilde{f}$  into the forward model results in a set of discrete values  $\tilde{g}_i$  that are consistent with the known measurements  $\gamma_i$ , once binarized. In the context of 1-bit compressed sensing, the enforcement of sign consistency has been originally proposed in [7], where a one-sided quadratic penalty function was considered. Trivial solutions were avoided by requiring that the signal lies on the unit sphere. Here, as in [8], we introduce a variational consistency principle that preserves the convexity of the problem without requiring additional non-convex constraints. Note that, although convexity is not required to ensure nontrivial solutions, it is exploited for the development of our algorithm and to ensure its convergence, as described in Section IV. Regarding the data-fidelity term, our contribution is to propose a penalty function  $\psi$  that is also suitable for bound optimization. We express our functional as

$$\mathcal{D}(\tilde{c}) = \sum_{i=1}^L \sum_{\mathbf{k}} \omega_i[\mathbf{k}] \psi(\tilde{g}_i[\mathbf{k}] \gamma_i[\mathbf{k}]), \quad (14)$$

where  $\tilde{g}_i$  and  $\tilde{c}$  are related in the same way as  $g_i$  and  $c$  in (10). The positive function  $\psi$  is defined as

$$\psi(t) = \begin{cases} M^{-1} - t, & t < 0 \\ M^{-1}(M^2 t^2 + Mt + 1)^{-1}, & \text{otherwise,} \end{cases} \quad (15)$$

where  $M$  is the total number of measurements. Besides penalizing sign inconsistencies, the rationale behind this definition is to yield nontrivial solutions while ensuring the convexity of the data term. The latter property holds because, according to (15), the Hessian of  $\mathcal{D}$  is well-defined and positive semidefinite [21]. The function  $\psi$  is itself  $\mathcal{C}^2$ -continuous and convex, its second derivative being always nonnegative. Moreover, this specific piecewise-rational polynomial function is suitable to the development of analytic upper bounds, as addressed in Section IV.

Given (4), negative arguments of  $\psi$  correspond to sign inconsistencies. As shown in Figure 3, our penalty function is linear in that regime. In that regard, the authors of [9] have shown that, in the binary compressed sensing framework, such an  $\ell_1$ -type penalty for consistency yields reconstructions that are of higher quality than with the  $\ell_2$  objective used in [7],

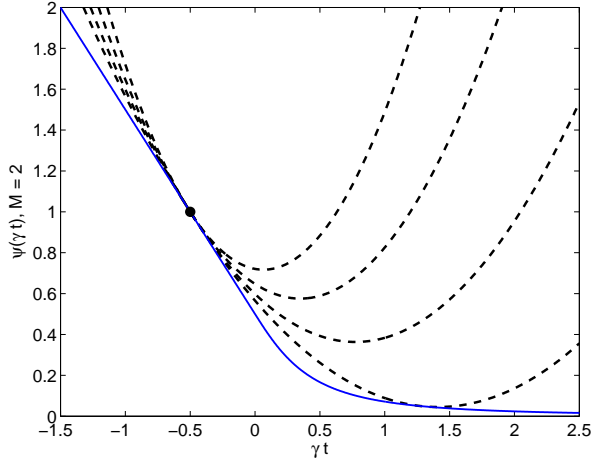


Fig. 3. Shape of our penalty function. As discussed in Section IV and further developed in Appendix A, the values  $\psi(\gamma t)$  (full line) can be bound from above by the quadratic function  $\psi_q(t|\tilde{g}^{(n)}, \gamma)$  around  $t = \tilde{g}^{(n)}$  (dot mark). Function values and derivatives must coincide at that point to satisfy (23). Among all possible parabolas (dashed lines), the solution  $\psi_q$  is the upper bound with infimum second derivative.

[18]. To some extent, these results confirm similar observations mentioned in [8]. This type of penalty also relates to the so-called *hinge loss* which is considered a better measure than the square loss for binary classification [9]. In our method, the values of the solution  $\tilde{c}$  are defined up to a common scale factor, and also up to an additive constant because  $\tau$  is not given. Non-constant solutions are favored by the contribution of the small nonlinear penalty that remains when the sign is correct. The transition between the linear and nonlinear regimes of  $\psi$  is  $\mathcal{C}^2$ -continuous and takes place at the origin. The applied penalty vanishes for increasingly positive arguments.

### B. Regularization Term

For inverse problems, it has been shown empirically that *frame-synthesis* regularization, which acts on transform-domain (e.g., wavelet) coefficients of the signal of interest, is outperformed by *frame-analysis* regularization, which directly operates on the signal itself [22], [23]. Accordingly, reconstruction algorithms often involve the latter approach when dealing with images; total-variation (TV) [24] is frequently used as a sparsifying transform [1], [25], [26]. Although suitable for regularization, the original form of TV is non-differentiable when the image gradient vanishes. As in the NESTA algorithm proposed in [27] for the recovery of sparse images, we therefore opt for a smooth approximation of the TV penalty based on a Huber potential function [28]. In order to guarantee the well-posedness of the problem, we also include an additional energy term in our expression, since the nullspace of  $\mathbf{A}$  can indeed be nonempty depending on  $\nabla_i$ . Approximating the Huber integral, our regularizer  $\mathcal{R}$  is then defined as

$$\mathcal{R}(\tilde{c}) = \sum_{\mathbf{k}} \mathcal{H}(\theta[\mathbf{k}]) + \lambda' \tilde{c}[\mathbf{k}]^2, \quad (16)$$

where each  $\theta[\mathbf{k}]$  is the norm of the gradient of  $\tilde{f}$  evaluated at position  $\mathbf{x} = \mathbf{k}$  and where  $\lambda'$  is a small positive constant. Based on a smoothing parameter  $\epsilon$ , the scaled Huber potential  $\mathcal{H}$  is defined as

$$\mathcal{H}(t) = \begin{cases} \epsilon^{-1}t^2, & |t| \leq \epsilon \\ 2|t| - \epsilon, & \text{otherwise.} \end{cases} \quad (17)$$

The gradient-norm sequence  $\theta$  is determined from the spatial derivatives  $\frac{\partial \tilde{f}}{\partial x_1}$  and  $\frac{\partial \tilde{f}}{\partial x_2}$  of the solution sampled in-between the grid nodes defined by the sum in (16). This type of discretization yields numerically stable solutions without oscillatory modes. It bears similarities with the so-called *marker-and-cell* methods used in fluid dynamics [29]. The expression of  $\theta$  as a function of  $\tilde{c}$  is

$$\theta[\mathbf{k}] = \sqrt{(\tilde{c} \star \beta_{x_1}^m)[\mathbf{k}]^2 + (\tilde{c} \star \beta_{x_2}^m)[\mathbf{k}]^2}, \quad (18)$$

where the  $\beta_{x_{1,2}}^m$  are directional B-spline-derivative filters defined as

$$\begin{aligned} \beta_{x_1}^m[\mathbf{k}] &= \beta^{m'}(k_1 + 1/2)\beta^m(k_2), \\ \beta_{x_2}^m[\mathbf{k}] &= \beta^{m'}(k_2 + 1/2)\beta^m(k_1). \end{aligned} \quad (19)$$

The first derivative  $\beta^{m'}$  of a B-spline has the symbolic expression given in [13].

## IV. RECONSTRUCTION ALGORITHM

### A. General Approach

In this section, we derive an algorithm to efficiently solve (13). Our main strategy is to recast the original formulation of the reconstruction problem as the partial minimization of successive quadratic costs  $\mathcal{J}_q$  that upper-bound  $\mathcal{J}$  locally around the current solution estimate  $\tilde{c}^{(n)}$ . Each  $\mathcal{J}_q$  can then be minimized using a specifically devised preconditioned conjugate-gradient method.

While sharing a common structure, every new quadratic cost is specified by the current solution. Its proper definition involves the pointwise nonlinear estimation of scalar quantities, which is a reweighting process akin to the one of iteratively reweighted least squares (IRLS). In our bound-optimization framework, each successive solution partially minimizes  $\mathcal{J}_q(\cdot|\tilde{c}^{(n)})$  with respect to its current value at  $\tilde{c}^{(n)}$ . Finding this solution amounts to partially solving a linear problem with a given initialization. We propose to precondition each of these linear problems according to its particular structure and find an approximate solution using the linear conjugate-gradient (CG) method. This approach ensures the global convergence of our method without having to specify any step parameter.

According to Figure 4, the successive reweighting and linear-resolution steps can be interpreted as alternate dequantization and deconvolution operations, respectively.

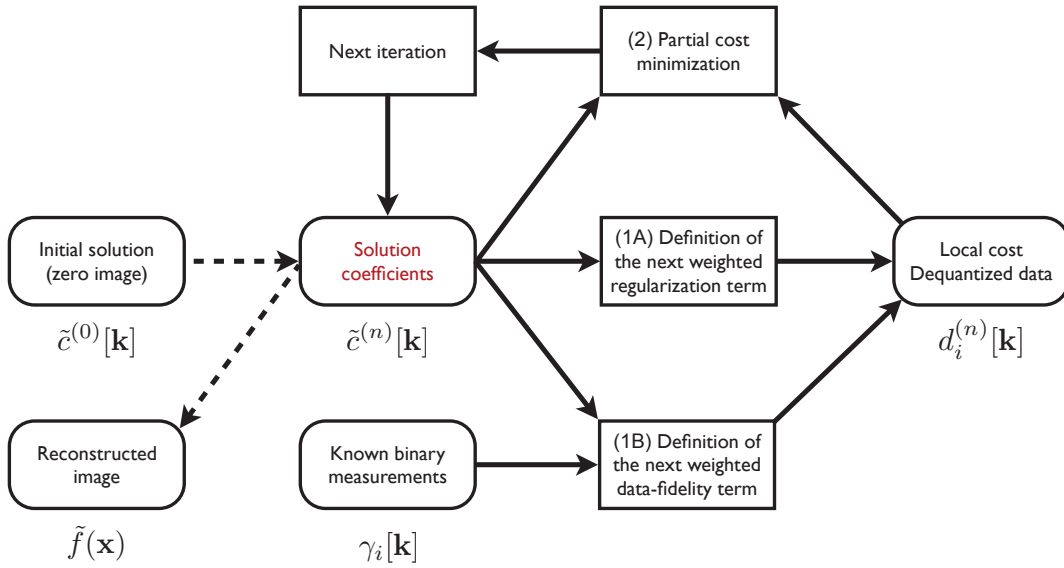


Fig. 4. Overall principle of our reconstruction algorithm. The solution coefficients are first initialized to zero and then updated by minimizing successive quadratic-cost functionals. Using the current solution  $\tilde{c}^{(n)}$ , Steps (1A) and (1B) determine the next local cost. Each of these two steps is related to a deconvolution problem where the data  $d_i$  to deconvolve correspond to dequantized versions of the available  $\gamma_i$ . An updated solution is found after minimization in Step (2). It determines the coefficients of the next solution. The overall convergence of the process is guaranteed because each quadratic cost is determined according to a bound-optimization approach and minimized using the current solution  $\tilde{c}^{(n)}$  as initialization.

### B. Upper Bound of the Data Term

In this part, we derive functionals of simpler form which upper-bound and approximate  $\mathcal{D}$  around some initial or current estimate of the solution. Following a *majorization-minimization* (MM) approach [30], we build the local quadratic cost  $\mathcal{D}_q^0(\cdot|\tilde{c}^{(n)})$  for the corresponding estimate  $\tilde{c}^{(n)}$  such that

$$\begin{aligned} \mathcal{D}_q^0(\tilde{c}^{(n)}|\tilde{c}^{(n)}) &= \mathcal{D}(\tilde{c}^{(n)}), \\ \mathcal{D}_q^0(\tilde{c}|\tilde{c}^{(n)}) &\geq \mathcal{D}(\tilde{c}). \end{aligned} \quad (20)$$

For convenience, we bound the cost by the penalty  $\psi$ . This fixes the structure of  $\mathcal{D}_q^0(\cdot|\tilde{c}^{(n)})$  as

$$\mathcal{D}_q^0(\tilde{c}|\tilde{c}^{(n)}) = \sum_{i=1}^L \sum_{\mathbf{k}} \omega_i[\mathbf{k}] \psi_q(\tilde{g}_i[\mathbf{k}]|\tilde{g}_i^{(n)}[\mathbf{k}], \gamma_i[\mathbf{k}]), \quad (21)$$

where  $\tilde{g}_i^{(n)}$  is the current estimate of  $\tilde{g}_i$  associated with the solution estimate  $\tilde{c}^{(n)}$ , and where  $\psi_q$  is a quadratic and scalar penalty function which takes the form

$$\psi_q(\tilde{g}_i|\tilde{g}_i^{(n)}, \gamma_i) = a_2(\tilde{g}_i^{(n)}, \gamma_i)\tilde{g}_i^2 + a_1(\tilde{g}_i^{(n)}, \gamma_i)\tilde{g}_i + a_0(\tilde{g}_i^{(n)}, \gamma_i), \quad (22)$$

where the  $a_j(\tilde{g}_i^{(n)}, \gamma_i)$  are polynomial coefficients. The values of  $\tilde{g}_i$  and  $\gamma_i$  depend on the solution estimate and the available binary measurements. Constraints (20) are then satisfied by fulfilling the simpler scalar conditions  $\forall \gamma \in \{-1, 1\}$  and  $\forall t \in \mathbb{R}$ ,

$$\begin{aligned} \psi_q(\tilde{g}_i|\tilde{g}_i^{(n)}, \gamma) &= \psi(\gamma\tilde{g}_i^{(n)}), \\ \psi_q(t|\tilde{g}_i^{(n)}, \gamma) &\geq \psi(\gamma t), \end{aligned} \quad (23)$$

where the subscripts have been dropped for convenience. These relations constrain the value of  $\psi_q$  and its derivative at  $\tilde{g}_i^{(n)}$ . As illustrated in Figure 3, further optimizing  $\psi_q$  to best approximate  $\psi(\gamma t)$  exhausts every remaining degree of freedom. This solution corresponds to the smallest positive  $a_2$  in (22) that allows (23) to be satisfied. The particular definition that we have proposed for the penalty function  $\psi$  allows for fast noniterative evaluation of the coefficients  $a_j$ . The actual expressions are derived in Appendix A. The resulting coefficients then specify the quadratic cost  $\mathcal{D}_q^0(\cdot|\tilde{c}^{(n)})$  as

$$\mathcal{D}_q^0(\tilde{c}|\tilde{c}^{(n)}) = \sum_{i=1}^L \sum_{\mathbf{k}} \omega_i[\mathbf{k}] w_i^{(n)}[\mathbf{k}] \left( \tilde{g}_i[\mathbf{k}] - d_i^{(n)}[\mathbf{k}] \right)^2 + \mathcal{K}, \quad (24)$$

where the scalar  $\mathcal{K}$  is constant with respect to  $\tilde{c}$ , and where  $w_i^{(n)}$  and  $d_i^{(n)}$  are sequences defined as

$$\begin{aligned} w_i^{(n)}[\mathbf{k}] &= a_2(\tilde{g}_i^{(n)}[\mathbf{k}], \gamma_i[\mathbf{k}]), \\ d_i^{(n)}[\mathbf{k}] &= -\frac{1}{2}(a_2^{-1}a_1)(\tilde{g}_i^{(n)}[\mathbf{k}], \gamma_i[\mathbf{k}]). \end{aligned} \quad (25)$$

Since the value of the constant  $\mathcal{K}$  is irrelevant for minimization, we define the cost  $\mathcal{D}_q(\cdot|\tilde{c}^{(n)})$  as  $\mathcal{D}_q^0(\cdot|\tilde{c}^{(n)})$  minus that constant. Dropping the subscript  $n$  for convenience, its explicit form in matrix notation as a function of the coefficients reduces to

$$\mathcal{D}_q(\tilde{c}|\tilde{c}^{(n)}) = \sum_{i=1}^L \left\| \mathbf{W}_i^{\frac{1}{2}} (\mathbf{A}_i \tilde{\mathbf{c}} - \mathbf{d}_i) \right\|_{\ell_2}^2, \quad (26)$$

where  $\mathbf{W}_i$  is a diagonal matrix with diagonal components  $\omega_i w_i^{(n)}$  and where  $\mathbf{d}_i$  is the vector associated with  $d_i^{(n)}$ .

### C. Upper Bound of the Regularizer

The Huber convex functional  $\mathcal{R}$  can be bound from above according to the same MM principles. The form of  $\mathcal{R}_q(\cdot|\tilde{\mathbf{c}}^{(n)})$  can be deduced from the results of [31]. Its matrix expression is

$$\mathcal{R}_q(\tilde{\mathbf{c}}|\tilde{\mathbf{c}}^{(n)}) = \lambda' \|\tilde{\mathbf{c}}\|_{\ell_2}^2 + \left\| \mathbf{W}_0^{\frac{1}{2}} \mathbf{R} \tilde{\mathbf{c}} \right\|_{\ell_2}^2, \quad (27)$$

where  $\mathbf{W}_0$  is a diagonal matrix with diagonal components

$$w_0^{(n)}[\mathbf{k}] = \max(\epsilon, \theta[\mathbf{k}])^{-1}, \quad (28)$$

and where  $\mathbf{R} = (\mathbf{R}_1, \mathbf{R}_2)$  is the discretized-gradient matrix. Each term  $\mathbf{R}_i$  is a circulant matrix associated with the filters  $\beta_{x_i}^m$  defined in (19).

### D. Quadratic-Cost Minimization

Combining the data and regularization terms (26) and (27), we obtain the local quadratic cost

$$\mathcal{J}_q(\tilde{\mathbf{c}}|\tilde{\mathbf{c}}^{(n)}) = \mathcal{D}_q(\tilde{\mathbf{c}}|\tilde{\mathbf{c}}^{(n)}) + \lambda \mathcal{R}_q(\tilde{\mathbf{c}}|\tilde{\mathbf{c}}^{(n)}). \quad (29)$$

In order to decrease  $\mathcal{J}$ , the new estimate  $\tilde{\mathbf{c}}^{(n+1)}$  must decrease  $\mathcal{J}_q(\cdot|\tilde{\mathbf{c}}^{(n)})$  itself. In other words, we have to satisfy

$$\mathcal{J}_q(\tilde{\mathbf{c}}|\tilde{\mathbf{c}}^{(n)}) \leq \mathcal{J}_q(\tilde{\mathbf{c}}^{(n)}|\tilde{\mathbf{c}}^{(n)}). \quad (30)$$

Defining  $\mathbf{I}' = \lambda \lambda' \mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix, the minimum of  $\mathcal{J}_q(\cdot|\tilde{\mathbf{c}}^{(n)})$  is the solution of

$$\mathbf{S} \tilde{\mathbf{c}} = \mathbf{y}, \quad (31)$$

with the system matrix

$$\mathbf{S} = \sum_{i=1}^L \mathbf{A}_i^T \mathbf{W}_i \mathbf{A}_i + \lambda \sum_{i=1}^2 \mathbf{R}_i^T \mathbf{W}_0 \mathbf{R}_i + \mathbf{I}' \quad (32)$$

and the right-hand-side vector

$$\mathbf{y} = \sum_{i=1}^L \mathbf{A}_i^T \mathbf{W}_i \mathbf{d}_i. \quad (33)$$

The huge matrix sizes entering into play require (31) to be solved iteratively. The positivity of  $w_i^{(n)}$  and  $w_0^{(n)}$  in (25) and (28) implies symmetry and positive-definiteness of  $\mathbf{S}$ , which allows for the CG method to be used. Initializing the latter at the current estimates, we guarantee the corresponding approximate solutions to comply with (30).

### E. Preconditioning

We also take advantage of preconditioning to obtain an approximate solution  $\tilde{\mathbf{c}}^{(n+1)}$  that is close to the exact minimum with fewer iterations. We impose our preconditioner  $\mathbf{P}$  to be a positive-definite circulant matrix, and define the two-sided preconditioned system

$$\mathbf{S}' = \mathbf{P}^{-\frac{1}{2}} \mathbf{S} \mathbf{P}^{-\frac{1}{2}}. \quad (34)$$

It is associated with the modified linear problem

$$\mathbf{S}' \tilde{\mathbf{c}}' = \mathbf{y}', \quad (35)$$

where  $\mathbf{y}'$  is predetermined as  $\mathbf{y}' = \mathbf{P}^{-\frac{1}{2}} \mathbf{y}$  and where the actual solution  $\tilde{\mathbf{c}}$  of the original problem is recovered as  $\tilde{\mathbf{c}} = \mathbf{P}^{-\frac{1}{2}} \tilde{\mathbf{c}}'$ . As a solution satisfying the above requirements, we consider

$$\mathbf{P} = \mathbf{F}^* \text{diag}(\mathbf{F} \mathbf{S} \mathbf{F}^*) \mathbf{F}, \quad (36)$$

where  $\mathbf{F}$  is the normalized DFT operator, where  $\mathbf{F}^*$  denotes its adjoint, and where  $\text{diag}(\cdot)$  is a projector onto the diagonal-matrix space. Definition (36) corresponds to the optimal circulant approximation of  $\mathbf{S}$  with respect to the Frobenius norm [32]. This solution is well-adapted to its convolutive nature as compared to diagonal preconditioning.

### F. Minimization Scheme

The successive quadratic bounds as well as the corresponding preconditioned linear problems being defined, we now describe the overall iterative minimization scheme that yields the solution  $\tilde{\mathbf{c}}$ , starting from an initialization  $\tilde{\mathbf{c}}^{(0)}$ . Our overall scheme is composed of two embedded iterative loops. The weight specification of the successive quadratic costs corresponds to external iterations with solutions  $\tilde{\mathbf{c}}^{(n)}$ .

Since our algorithm involves upper bounds that are partially minimized and that satisfy MM conditions of the form (20), it is part of the generalized MM (GMM) family [30]. In that regard, the continuity of our functional  $\mathcal{J}_q$  implies that the MM sequence  $\{\mathcal{J}(\tilde{\mathbf{c}}^{(0)}), \mathcal{J}(\tilde{\mathbf{c}}^{(1)}), \mathcal{J}(\tilde{\mathbf{c}}^{(2)}), \dots\}$  converges monotonically to a stationary point of  $\mathcal{J}$ . The convexity of  $\mathcal{J}$  also implies that the whole minimization process is compatible with Nesterov's acceleration technique [20], which we apply to update our estimates. This requires the use of auxiliary solutions that we mark with star subscripts, as well as the definition of scalar values  $\sigma^{(n)}$ . The steps of our global scheme yielding the solution  $\tilde{\mathbf{c}}$  are described in Figure 5.

We use  $\mathcal{I}_{ext}$  external iterations, each of which corresponds to a refined quadratic approximation  $\mathcal{J}_q$  of the global convex cost. For the partial resolution of each internal problem, we apply CG on the modified system (35). Accordingly, the corresponding intermediate value  $\tilde{\mathbf{c}}'$  is first initialized to the current solution estimate in the preconditioned domain, and then updated using  $\mathcal{I}_{int}$  CG iterations each time. In accordance with (9), the final continuous-domain image is obtained from the coefficients  $\tilde{\mathbf{c}}$  as

$$\tilde{f}(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{Z}^2} \tilde{c}[\mathbf{k}] \beta^m(\mathbf{x} - \mathbf{k}). \quad (37)$$

As demonstrated in Section V-A, the use of Nesterov's technique and of preconditioning to solve the linear problems ensure the fast convergence of our method.

## V. EXPERIMENTS

We conduct experiments on grayscale images that are part of a standard test set. First, we evaluate the computational performance of our algorithm in Section V-A and show baseline results in Section V-B. In Section V-C, we propose



```

1) Initial  $\tilde{\mathbf{c}}^{(0)}$  taken as the zero vector
2) Initial values  $\tilde{\mathbf{c}}_*^{(0)} = \tilde{\mathbf{c}}^{(0)}$ ,  $n = 0$ ,  $\sigma^{(0)} = 1$ 
while  $n < \mathcal{I}_{ext}$  do
  a) Specification of  $\mathbf{S}$  and  $\mathbf{y}$  given  $\tilde{\mathbf{c}}^{(n)}$ 
  b) Computation of the preconditioner  $\mathbf{P}$  linked to  $\mathbf{S}'$ 
  c) External count  $n \leftarrow n + 1$ 
  d) Computation of  $\mathbf{y}' = \mathbf{P}^{-\frac{1}{2}}\mathbf{y}$ 
  e) Internal initialization  $\tilde{\mathbf{c}}' = \mathbf{P}^{\frac{1}{2}}\tilde{\mathbf{c}}^{(n-1)}$ 
  f) Update of  $\tilde{\mathbf{c}}'$  with  $\mathcal{I}_{int}$  CG iterations on  $\mathbf{S}'\tilde{\mathbf{c}}' = \mathbf{y}'$ 
  g) Nesterov's solution update  $\tilde{\mathbf{c}}_*^{(n)} = \mathbf{P}^{-\frac{1}{2}}\tilde{\mathbf{c}}'$ 
  h) Nesterov's step update

$$\sigma^{(n)} = \frac{1}{2} + \sqrt{\frac{1}{4} + (\sigma^{(n-1)})^2}$$

  i)  $\tilde{\mathbf{c}}_* = \tilde{\mathbf{c}}_*^{(n)} - \tilde{\mathbf{c}}_*^{(n-1)}$ 
  j)  $\tilde{\mathbf{c}}^{(n)} = \tilde{\mathbf{c}}_*^{(n)} + \sigma^{(n)-1}(\sigma^{(n-1)} - 1)\tilde{\mathbf{c}}_*$ 
end
3) Solution coefficients  $\tilde{\mathbf{c}} = \tilde{\mathbf{c}}^{(n)}$ 

```

Fig. 5. Minimization approach described in matrix notation.

an estimate of the acquisition quality based on the spatial redundancy of the available measurements. In Sections V-D and V-E, we address cases where downsampling and finite differentiation are used for data acquisition. In particular, we determine to what extent these strategies impact on the acquisition and reconstruction quality. We finally assess the optimal rate-distortion performance of our method for distinct amounts of measurements in Section V-F.

The discretization (9) does not induce any loss because we match the square grid of  $N_0 \times N_0$  spline coefficients to the resolution of each digital test image, choosing  $m = 1$ . Specifically, we determine  $c$  beforehand such that  $f$  interpolates the corresponding pixel values<sup>4</sup>. In order to maximize the acquisition bandwidth, the size  $K \times K$  of the phase mask and the number  $M_0 \times M_0$  of sensors are themselves set to  $N_0 \times N_0$ . The sampling prefilter  $\phi$  is defined as a 2D separable rectangular window. The threshold  $\tau$  is set to the mean image intensity<sup>5</sup> when no finite differentiation is used, and to zero otherwise. The latter choice is a heuristic that directly yields equidistributed binary measurements  $\gamma_i$  from our data as in [8], without requiring any optimization or further refinement. For non-unit  $\Lambda$ , we consider identical spatial masks  $\omega_i$  that correspond to horizontal and vertical subsampling, which allows for the proper display and evaluation of our measurements. Our reconstruction parameters are  $\lambda = 10^{-4}$ ,  $\lambda' = 10^{-5}$ ,  $\epsilon = 5 \cdot 10^{-4}$ ,  $\mathcal{I}_{ext} = 20$ , and  $\mathcal{I}_{int} = 4$ . The smoothing parameter  $\epsilon$  chosen for our regularizer aims at approximating TV as in [27], while the small values of the constants  $\lambda$  and  $\lambda'$  ensure that the reconstructions are consistent with the binary measurements with enough accuracy (*i.e.*, about 99% or above).

We have found that the most-consistent solutions are also the ones of highest quality, which corroborates the results of

<sup>4</sup>Given our forward model and the high values of  $N_0$  involved in our experiments, the choice of  $m$  has no significant impact.

<sup>5</sup>This quantity corresponds to the mean component value of the vector  $\mathbf{g}$ . It is assumed to be known for reconstruction.

[9]. Knowing that each instance of (35) can be solved partially, the choice of  $\mathcal{I}_{int}$  is meant to maximize computational performance, while the value of  $\mathcal{I}_{ext}$  is used as a stop criterion. Note that the values of  $\epsilon$  and  $\lambda$  cannot be reduced further without impacting negatively on the speed of convergence.

In order to provide a quality assessment in terms of signal-to-noise ratio (SNR), the mean and variance of the solution coefficients are matched to the reference signal. We also define a quantity called blockwise-corrected SNR (BSNR) where this same matching is performed blockwise using  $8 \times 8$  blocks. As discussed in Section V-D, the BSNR is consistent with visual perception.

### A. Computational Performance

To evaluate the computational performance of our algorithm, we perform a reconstruction experiment on a  $256 \times 256$  test image using  $M_0^2 = 256^2$ ,  $L = 1$ ,  $\Lambda = 1$ , and no finite differentiation. The results are reported in Figure 6, including a comparison with the BIHT algorithm<sup>6</sup> introduced for reconstruction from binary measurements in [9]. These results demonstrate that Nesterov's acceleration method, as well as the preconditioning used in our algorithm, play a central role to obtain fast convergence. By contrast, we have observed that 3,000 iterations are required to ensure convergence with BIHT—which is used for the experiments of Section V-D—as opposed to a total of  $\mathcal{I}_{ext}\mathcal{I}_{int} = 80$  internal iterations with our algorithm. This corresponds to an order-of-magnitude improvement in time efficiency.

### B. Baseline Results

Our framework can handle several measurement sequences unlike in [8]. Accordingly, the goal in this part is to reconstruct the  $512 \times 512$  images *Lena* and *Barbara* from distinct numbers  $L$  of acquisitions with  $\Lambda = 1$  and no finite differentiation. Each acquisition includes  $M_0^2 = 512^2$  samples, the total number of measurements being multiplied by the corresponding  $L$ .

The binary acquisitions and the corresponding reconstructions with our algorithm are shown in the spatial domain in Figure 7. In both examples, the reconstruction quality substantially improves with  $L$ , one single acquisition being already sufficient to preserve substantial grayscale and edge information. The binary measurements of Figure 7 are not interpretable visually because the image information has been spread out through the filters  $h_i$ . These measurements follow a random distribution that originates from the pseudo-random phases  $\nu_i$  of the masks, and that is heavily correlated spatially as in [8]. As a matter of fact, random-convolution measurements do not display strict statistical incoherence [1]. We investigate below how spatial correlation can be quantified and reduced to improve reconstruction.

<sup>6</sup>We have adapted BIHT to our forward model, assuming sparsity in the Haar-wavelet domain. Besides its simplicity, the latter choice was observed to yield higher-quality results in our case than when using higher-order Daubechies wavelets, despite the generated block artifacts. Each iteration involves a gradient step scaled as  $M^{-1/2}\|\mathbf{A}\|_2^{-1}$  and renormalization [9]. A zero-mean  $\mathbf{A}$  is used in the algorithm to handle the case where  $\tau$  is nonzero. The sparsity-level parameter specifying the assumed amount of nonzero wavelet coefficients is set as 2,000. Both BIHT and the proposed algorithms have been implemented in MATLAB.



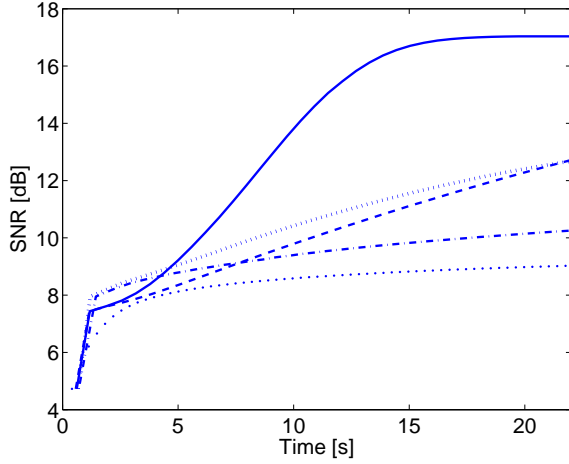


Fig. 6. Reconstruction SNR as a function of time for *Montage* ( $256 \times 256$ ). For our reconstruction method, the sole use of preconditioning (dashed line) or Nesterov’s acceleration (dotted line) already improves the convergence rates as compared to standard CG (mixed line). When both techniques are enabled (solid line), the performance of our algorithm improves substantially. For comparison, the reconstruction performance of BIHT is also shown for the same problem (bottom dots). In the latter case, each corresponding iteration lasts about half a second. The times that are given correspond to an execution of the algorithms on Mac OS X version 10.7.1 (MATLAB R2011b) with a Quad-Core Xeon  $2 \times 2.8$  GHz and 4 GB of DDR2 memory.

### C. Incoherence Estimation

The potential quality of reconstruction depends on the appropriateness of  $\mathbf{A}$  for binary compressed sensing. We assume our matrix to be suitable for the specific data in hand when the corresponding binarized measurements behave as independent and identically distributed random variables. As a practical solution, we propose to estimate the “randomness” of the acquired  $\gamma_i$  through their autocorrelation [33]. We specifically infer a correlation distance  $\alpha$  based on the unnormalized autocorrelations  $\rho_i^*$  of our (possibly subsampled) binary sequences  $\gamma_i^*$ . This distance is used as a quality indicator, inasmuch as it measures the degree of spatial redundancy arising in our measurements. To determine this value, we first compute the characteristic length  $\alpha^i$  of each autocorrelation peak, using the standard deviation of  $|\rho_i^*|^4$  for the sake of robustness. The autocorrelation being symmetric and centered at the origin, we write that

$$\alpha^i = \left( \frac{\sum_{\mathbf{k}} |\rho_i^*[\mathbf{k}]|^4 \|\mathbf{k}\|^2}{\sum_{\mathbf{k}} |\rho_i^*[\mathbf{k}]|^4} \right)^{1/2}. \quad (38)$$

Averaging  $\alpha^i$  over  $i$  then yields the final  $\alpha$ . As shown in the sequel, this value strongly depends on the parameters of the forward model. In particular, it can be decreased compared to the case of Section V-B by enabling downsampling (*i.e.*, non-unit  $\Lambda$ ) or finite differentiation in our framework. Note that, as in [8], our choice for the threshold  $\tau$  ensures the uniformity of the binary distribution of the measurements.

### D. Influence of Acquisition Modality

In this section, we investigate the performance of finite differentiation when used in our framework. To this end, we

choose a fixed set of two perpendicular first-derivative filters whose  $\mathcal{Z}$ -domain expressions are  $(z_1 - z_1^{-1})$  for the horizontal orientation and  $(z_2 - z_2^{-1})$  for the vertical orientation, respectively. Assuming an even  $L$ , the former filter is applied on acquisition sequences of even index, and the latter one is applied on the remaining indices. The operation of each filter  $\nabla_i$  followed by zero thresholding is physically realizable by means of binary comparators that are connected to the two corresponding pixels. From a practical standpoint, such an approach eliminates the need of threshold calibration.

In order to compare the acquisition modalities with and without finite differentiation, we perform experiments on several  $256 \times 256$  images. These experiments involve  $M = 131,072$  measurements taken in  $L = 2$  acquisitions, using  $M_0^2 = 256^2$  and  $\Lambda = 1$ . Besides our own algorithm, BIHT is also considered for reconstruction in each case. The results are reported in Table I, and shown in Figure 8 for *House*. The best numerical values are emphasized in the tables using bold notation.

Our qualitative and quantitative results demonstrate that finite differentiation globally yields the best reconstructions. These solutions consistently correspond to lower  $\alpha$  values as well, which reflects itself visually in less-redundant binary measurements. Finite differentiation decreases redundancy because it spatially decorrelates the image measurements  $g_i$  before quantization. Because finite differentiation senses the high-frequency content of the measurements, most visual features such as edges are indeed better restored as compared to the other acquisition modalities. In return, reconstructions tend to display slightly higher low-frequency error. Because of its cumulative nature, the latter may then cause substantial SNR deterioration in unfavorable scenarios. In such cases, however, the amount of visual details is still higher, as illustrated in Figure 8. For instance, fine details such as the house gutter are better preserved. We observe that the BSNR measure is consistent with visual impression, as it adapts to slow intensity drifts in the solution. For both acquisition modalities, our algorithm based on TV yields the best reconstructions. This confirms the suitability of TV for our problem, in accordance with the discussion of Section III-B. Note, however, that proper adjustment of the sparsity level in BIHT is delicate. For instance, images that are sparser than the assumed level might lead to suboptimal reconstructions in Table I.

### E. Respective Influence of $\Lambda$ and $L$

The following experiments address how reconstruction quality can be maximized given a fixed measurement budget, using the same  $256 \times 256$  images as above. Considering the finite-differentiation modality specified in Section V-D, our strategy is to further decrease spatial redundancy by sharing the measurements between more acquisitions. Choosing  $M_0^2 = 256^2$  and  $M = 32,768$  as constraints, we thus adapt the ratio  $\Lambda$  to the number of acquisitions as  $\Lambda^{-1} = 2L$ . On the one hand, minimizing  $L$  reduces to previous system configurations. On the other hand, maximizing it is highly inefficient, as it amounts to taking one single measurement per convolutive acquisition. A tradeoff has to be found between these two



Fig. 7. Results on *Lena* and *Barbara* ( $512 \times 512$ ) for distinct numbers  $L$  of acquisitions using  $M_0^2 = 512^2$  and  $\Lambda = 1$  without finite differentiation ( $M = L \cdot 512^2$  measurements in total). First row, from left to right: first acquisition  $\gamma_1$  of *Lena* using our model, and reconstruction from one ( $M = 262,144$ , SNR: 17.49 dB, BSNR: 22.35 dB), two ( $M = 524,288$ , SNR: 22.42 dB, BSNR: 24.61 dB), and four ( $M = 1,048,576$ , SNR: 26.46 dB, BSNR: 27.13 dB) acquisitions. Second row: first acquisition  $\gamma_1$  of *Barbara* using our model, and reconstruction from one ( $M = 262,144$ , SNR: 13.96 dB, BSNR: 16.09 dB), two ( $M = 524,288$ , SNR: 17.69 dB, BSNR: 17.74 dB), and four ( $M = 1,048,576$ , SNR: 20.3 dB, BSNR: 20.28 dB) acquisitions.

Modality Reconstruction	Standard Approach					Finite Differences				
	Proposed (TV)		BIHT (Haar)		$\alpha$	Proposed (TV)		BIHT (Haar)		$\alpha$
SNR	BSNR	SNR	BSNR	SNR		BSNR	SNR	BSNR		
<i>Bird</i>	25.64	27.80	19.80	22.56	54.00	<b>25.81</b>	<b>31.66</b>	15.17	23.88	<b>33.08</b>
<i>Cameraman</i>	20.65	20.96	15.95	16.32	64.99	<b>22.63</b>	<b>24.04</b>	5.87	17.16	<b>16.04</b>
<i>House</i>	<b>25.67</b>	26.44	20.40	21.58	47.82	24.38	<b>28.85</b>	13.83	22.30	<b>20.16</b>
<i>Peppers</i>	<b>20.16</b>	21.79	14.71	15.43	40.30	18.21	<b>24.95</b>	7.15	15.61	<b>19.87</b>
<i>Shepp-Logan</i>	19.25	20.00	9.53	9.95	34.15	<b>22.96</b>	<b>25.24</b>	5.72	12.26	<b>11.58</b>

TABLE I  
ACQUISITION MODALITIES COMPARED ON  $256 \times 256$  IMAGES USING  $M_0^2 = 256^2$ ,  $L = 2$ , AND  $\Lambda = 1$  ( $M = 131,072$ ).

limits to improve the quality of the reconstructions while preserving the parallelism of our model.

Our numerical results are reported in Table II, the measurements and reconstruction of *Peppers* being shown for two distinct settings in Figure 9. The values of Table II confirm that the correlation length  $\alpha$  consistently decreases with  $\Lambda$ . Moreover, the SNR and BSNR improve by several decibels when increasing  $L$ . This is further corroborated by the visual results of Figure 9. In particular, grayscale information is more-finely preserved in the solution displayed on the right. Interestingly, the increase in quality starts saturating when  $\alpha$  reaches near-optimal values, as shown in Table II. The compression performance of our method is thus optimal or nearly optimal with  $L \geq 8$  for a given amount of measurements. These results confirm the strong inverse correlation between measurement redundancy and reconstruction quality.

#### F. Rate-Distortion Performance

In this section, we confront our global acquisition and reconstruction framework (GF) with the single-convolution

framework (SF) of [8]. The following experiments allow us to evaluate their respective image-reconstruction performance in terms of the rate of distortion, defining the number of bits per pixel (bpp) as the ratio between  $M$  and the raw bitsize of the corresponding uncompressed 8-bit-grayscale image.

In order to decrease  $\alpha$  within a reasonable amount of acquisitions, our forward model is parameterized with  $L = 8$  and  $\Lambda = L^{-1}x$ , depending on the chosen bitrate  $x$  in bpp. The number of measurements taken on an  $N_0 \times N_0$  test image is thus  $M = xN_0^2$  since  $M_0 = N_0$ . Our method is evaluated with (D) and without (S) finite differentiation. In the SF case, the sensor resolution has to match  $M$  strictly, because one single convolution is performed without subsequent drop of samples. The forward model is configured accordingly, adapting the remaining parameters to the image size as in our method. That particular framework requires equal rational factors for resampling, which implies that certain bitrates cannot be evaluated. The reconstruction parameters are set as in the last experiment of [8].

Results on several test images are reported in Table III.

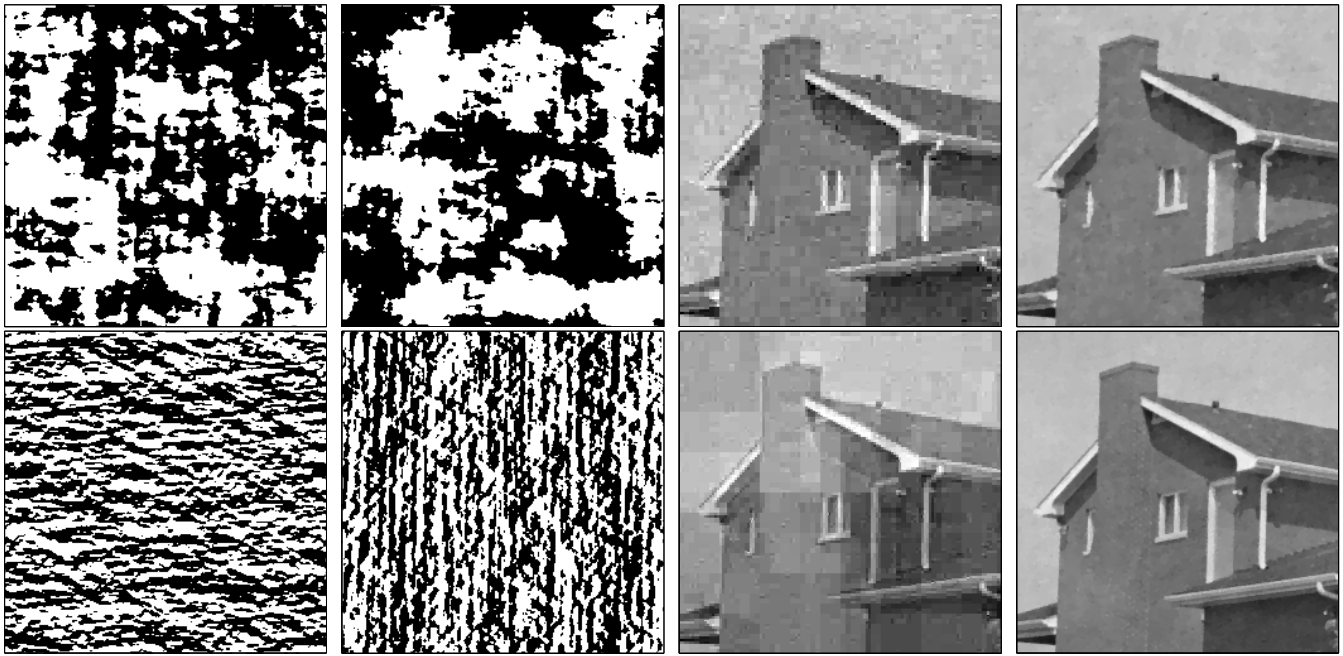


Fig. 8. Acquisition modalities compared on *House* ( $256 \times 256$ ) using  $M_0^2 = 256^2$ ,  $L = 2$ , and  $\Lambda = 1$  ( $M = 131,072$ ). First row, from left to right: acquisitions  $\gamma_i$  without finite differentiation, and reconstruction using BIHT (SNR: 20.40 dB, BSNR: 21.58 dB) and our algorithm (SNR: 25.67 dB, BSNR: 26.44 dB). Second row: acquisitions  $\gamma_i$  with finite differentiation, and reconstruction using BIHT (SNR: 13.83 dB, BSNR: 22.3 dB) and our algorithm (SNR: 24.38 dB, BSNR: 28.85 dB).

Parameters	$L = 2, \Lambda = 1/4$			$L = 4, \Lambda = 1/8$			$L = 8, \Lambda = 1/16$			$L = 16, \Lambda = 1/32$			$L = 32, \Lambda = 1/64$		
	SNR	BSNR	$\alpha$	SNR	BSNR	$\alpha$	SNR	BSNR	$\alpha$	SNR	BSNR	$\alpha$	SNR	BSNR	$\alpha$
<i>Bird</i>	22.62	28.89	13.76	22.77	29.25	10.57	24.30	29.41	5.66	25.35	<b>29.57</b>	4.94	<b>25.37</b>	29.49	<b>2.31</b>
<i>Cameraman</i>	18.73	20.79	5.91	18.63	21.08	3.77	<b>19.91</b>	21.30	1.92	19.81	21.26	1.59	19.53	<b>21.38</b>	<b>1.04</b>
<i>House</i>	20.71	26.34	8.05	21.10	26.51	6.35	24.01	26.81	3.74	24.05	26.88	2.83	<b>24.56</b>	<b>26.96</b>	<b>1.78</b>
<i>Peppers</i>	15.09	21.29	8.01	15.68	21.98	6.14	18.95	22.28	2.99	19.01	22.42	2.63	<b>19.19</b>	<b>22.47</b>	<b>1.49</b>
<i>Shepp-Logan</i>	16.88	19.42	4.26	16.84	19.50	2.59	17.20	19.60	1.51	17.48	<b>19.64</b>	1.27	<b>17.49</b>	19.58	<b>0.93</b>

TABLE II

INFLUENCE OF  $\Lambda$  AND  $L$  EVALUATED ON  $256 \times 256$  IMAGES USING  $M_0^2 = 256^2$  AND FINITE DIFFERENTIATION. THE SAME NUMBER OF MEASUREMENTS  $M = 32,768$  IS SHARED BETWEEN DISTINCT NUMBERS OF ACQUISITIONS ( $M/N = 1/2$ ).

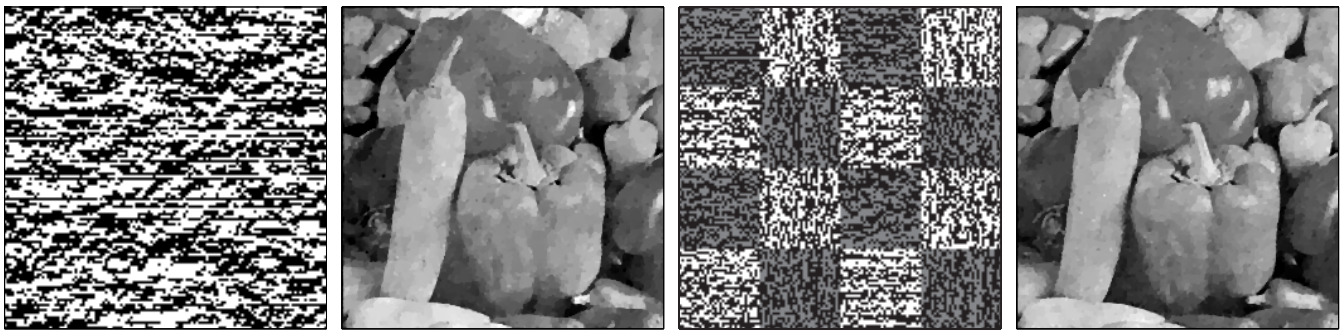


Fig. 9. Results on *Peppers* ( $256 \times 256$ ) when sharing  $M = 32,768$  measurements between distinct numbers of acquisitions with finite differentiation and  $M_0^2 = 256^2$ . From left to right: acquisition and reconstruction for  $L = 2$  and  $\Lambda = 1/4$  with  $\gamma_1$  (SNR: 15.09 dB, BSNR: 21.29 dB), and for  $L = 32$  and  $\Lambda = 1/64$  with  $\gamma_1$  to  $\gamma_{16}$  shown in concatenated form using a gray/white checkerboard-type display (SNR: 19.19 dB, BSNR: 22.47 dB).

They indicate that at least one version of our method always exceeds SF in terms of reconstruction quality. This confirms the relevance of sharing the acquired data between more acquisitions as a means of decreasing spatial redundancy. This strategy thus tends to compensate the non-ideal statistical properties of binary measurements that are based on random convolutions. In the case of SF, spatial redundancy cannot be

decreased similarly since only one convolutive acquisition is used. As previously observed, the (D) modality of our method can yield worse SNR values in certain configurations, while displaying superior BSNR performance globally. Nevertheless, these complementary results reveal an advantageous SNR performance of (D) at higher bitrates.

The efficiency of our method at  $1/8$  bpp, which corresponds



to a compression factor of 64, is illustrated for both modalities in Figure 10. Also shown is the plain JPEG version of the image compressed at similar bitrate. In this example, the GF framework with finite differentiation yields the best BSNR. We observe that the corresponding reconstruction contains fine details despite the low amount of measurements. It is also visually more pleasant than the JPEG solution. This experiment illustrates the highest compression ratio at which our method reconstructs images with reasonable quality. From a general standpoint, the results of this section demonstrate that, although generally inferior, the rate-distortion performance of binary compressed sensing can compete with JPEG at low bitrates. This can be deduced by comparing the plain-JPEG performance to the corresponding SNR values reported in Table III and corroborates the analysis of [34] where compressed sensing is compared to traditional image-compression methods.

## VI. CONCLUSIONS

We have proposed a binary compressed-sensing framework which is suitable for images. In our experiments, we have illustrated how measurement redundancy can be minimized by properly configuring our acquisition model. We have considered the single-acquisition case as well as a multi-acquisition strategy. In the two cases, our reconstruction algorithm has demonstrated state-of-the-art reconstruction performance on standard images. In particular, detailed features have been successfully recovered from small amounts of binary data. From a global perspective, our results confirm the 1-bit-compressed-sensing paradigm to be promising for imaging applications. In that regard, the specific interest of our method is to involve binary measurements that are suitable to convex optimization. We have proposed an iterative algorithm that combines preconditioning and Nesterov's approach to provide very efficient reconstructions of our measurements. Synthetic experiments demonstrate the potential of our method.

## ACKNOWLEDGMENTS

The authors are most indebted to Christian Bovet, Jean-Paul Cano, and Anthony Saugey (Essilor International, France) for their support and guidance. They also thank Christophe Moser (Ecole polytechnique fédérale de Lausanne, Switzerland) and Chandra Sekhar Seelamantula (Indian Institute of Science, Bangalore) for fruitful discussions.

## APPENDIX A

### COEFFICIENTS OF THE PENALTY BOUNDS

#### A. Formulation of the Optimization Task

The continuity of  $\psi(\gamma t)$  and the upper-bound conditions on  $\psi_q(t|\tilde{g}^{(n)}, \gamma)$  impose that the value and first derivative of these two functions coincide at  $t = \tilde{g}^{(n)}$ . This requires that

$$\begin{aligned} a_0 &= \psi(\gamma\tilde{g}^{(n)}) - \tilde{g}^{(n)}(\tilde{g}^{(n)}a_2 + a_1), \\ a_1 &= -2\tilde{g}^{(n)}a_2 + \psi'(\gamma\tilde{g}^{(n)}). \end{aligned} \quad (39)$$

The remaining degree of freedom  $a_2 \in \mathbb{R}_+ \setminus \{0\}$  is optimized so as to best approximate  $\psi$ . The resulting optimal  $a_2$  corresponds to the lowest positive value satisfying (23). In that configuration, the parabola  $\psi_q(t|\tilde{g}^{(n)}, \gamma)$  touches one and only one distinct point of  $\psi(\gamma t)$  at  $t = \tilde{g}^T$ . The convexity of  $\psi$  ensures the existence and uniqueness of the solution.

#### B. Solution

According to (22), the abscissas of the intersections between  $\psi$  and  $\psi_q(\cdot|\tilde{g}^{(n)}, \gamma)$  are solutions of

$$a_2t^2 + a_1t + a_0 = \psi(\gamma t). \quad (40)$$

These solutions correspond to the set union

$$\mathcal{S} = \{t \leq 0 : \mathcal{P}_1(t) = 0\} \cup \{t > 0 : \mathcal{P}_2(t) = 0\}, \quad (41)$$

where  $\mathcal{P}_{1,2}(t) = 0$  gives the intersections between  $\psi_q(\cdot|\tilde{g}^{(n)}, \gamma)$  and the linear and nonlinear parts of  $\psi$ . This corresponds to the separate formulas of (15) without the argument condition. Accordingly, the polynomials  $\mathcal{P}_{1,2}$  are expressed as

$$\begin{aligned} \mathcal{P}_1(t) &= a_2t^2 + (a_1 + \gamma)t + (a_0 - M^{-1}), \\ \mathcal{P}_2(t) &= (M^2t^2 + M\gamma t + 1)(a_2t^2 + a_1t + a_0) - M^{-1}. \end{aligned} \quad (42)$$

The optimal  $\psi_q(t|\tilde{g}^{(n)}, \gamma)$  is tangent to  $\psi(\gamma t)$  at  $t = \tilde{g}^{(n)}, \tilde{g}^T \in \mathcal{S}$ , and intersects no other point. This causes the two double roots  $\tilde{g}^{(n)}$  and  $\tilde{g}^T$  to appear in one of the two polynomials, be it jointly or not. Either of these two roots cancels the discriminant  $D$  of the associated polynomial. For the sake of conciseness, we define  $u = M\gamma\tilde{g}^{(n)}$  and consider two distinct cases.

1) *The point  $\tilde{g}^{(n)}$  is in the nonlinear part of  $\psi$ :* In this case, where  $u \geq 0$ , the coefficients  $a_0$  and  $a_1$  are expressed as

$$\begin{aligned} a_0 &= M^{-1}((u^2 + u + 1)^{-1} - M^{-1}u^2a_2 - \gamma ua_1), \\ a_1 &= -\gamma(2M^{-1}ua_2 + (2u + 1)(u^2 + u + 1)^{-2}). \end{aligned} \quad (43)$$

Then, the optimal parabola can be tangent at a distinct point of  $\psi$  either in the same nonlinear part, or in the linear part. If  $\tilde{g}^T$  lies in the linear part, the corresponding polynomial  $\mathcal{P}_1$  contains one double root  $\tilde{g}^T$  for an optimal  $a_2$ . This first subcase corresponds to the solution

$$\begin{aligned} a_2' &= \{a_2 \in \mathbb{R}_+^* : D(\mathcal{P}_1(\cdot)) = 0\} \\ &= \frac{1}{4}M \frac{u(u^2 + 2u + 3)^2}{(u^2 + u + 1)^3}. \end{aligned} \quad (44)$$

If  $\tilde{g}^T$  lies in the nonlinear part of  $\psi$ , the corresponding  $\mathcal{P}_2$  contains two double roots. Its discriminant is thus always zero regardless of  $a_2$ . Nevertheless, this same quantity divided by  $(t - \tilde{g}^{(n)})^2$  is a viable indicator, as it only vanishes in the optimal case. This yields the solution for this second subcase as



Bitrate Sampling Ratio $\Lambda^*$		1/16 bpp 1/128	1/8 bpp 1/64	1/4 bpp 1/32	1/2 bpp 1/16	1 bpp 1/8	2 bpp 1/4	4 bpp 1/2
Image	Method	SNR / BSNR						
<i>Bird</i> (256 × 256)	SF	18.35 / 21.85	-	21.27 / 24.44	-	22.95 / 26.05	-	23.78 / 27.14
	GF (S)	<b>19.44</b> / 21.94	<b>21.65</b> / 23.40	<b>23.74</b> / 25.04	<b>25.58</b> / 26.41	27.13 / 27.76	28.36 / 28.96	28.58 / 29.57
	GF (D)	16.84 / <b>23.56</b>	19.79 / <b>25.65</b>	22.30 / <b>27.65</b>	24.30 / <b>29.41</b>	<b>27.34</b> / <b>31.19</b>	<b>30.67</b> / <b>33.07</b>	<b>33.17</b> / <b>34.54</b>
<i>Cameraman</i> (256 × 256)	SF	13.86 / 14.79	-	16.34 / 17.02	-	18.14 / 18.94	-	19.33 / 20.34
	GF (S)	<b>14.63</b> / 15.03	<b>16.06</b> / 16.03	<b>17.26</b> / 17.22	18.54 / 18.42	19.78 / 19.68	21.27 / 21.22	22.67 / 22.73
	GF (D)	11.33 / <b>16.11</b>	15.00 / <b>17.84</b>	17.20 / <b>19.57</b>	<b>19.91</b> / <b>21.30</b>	<b>21.96</b> / <b>23.01</b>	<b>23.73</b> / <b>24.63</b>	<b>25.72</b> / <b>26.09</b>
<i>House</i> (256 × 256)	SF	17.36 / 20.39	-	20.74 / 23.08	-	23.21 / 25.10	-	24.60 / 26.31
	GF (S)	<b>18.39</b> / 20.40	<b>20.62</b> / 21.87	<b>22.67</b> / 23.32	<b>24.47</b> / 24.73	25.74 / 25.85	27.11 / 27.11	27.78 / 27.87
	GF (D)	15.48 / <b>21.25</b>	18.56 / <b>23.47</b>	20.94 / <b>25.15</b>	24.01 / <b>26.81</b>	<b>26.62</b> / <b>28.23</b>	<b>28.78</b> / <b>29.80</b>	<b>30.39</b> / <b>31.10</b>
<i>Peppers</i> (256 × 256)	SF	12.43 / 14.63	-	15.09 / 16.94	-	17.35 / 19.85	-	18.58 / 21.49
	GF (S)	<b>14.31</b> / 15.38	<b>15.99</b> / 16.66	<b>17.44</b> / 17.80	<b>19.02</b> / 19.38	20.91 / 21.36	23.06 / 23.51	24.67 / 25.02
	GF (D)	10.84 / <b>15.70</b>	13.40 / <b>17.80</b>	16.11 / <b>19.97</b>	18.95 / <b>22.28</b>	<b>21.20</b> / <b>24.57</b>	<b>23.87</b> / <b>26.58</b>	<b>26.73</b> / <b>28.26</b>
<i>Shepp-Logan</i> (256 × 256)	SF	7.28 / 8.73	-	12.57 / 13.88	-	17.33 / 18.24	-	22.33 / 22.89
	GF (S)	<b>8.52</b> / 10.21	<b>10.95</b> / 12.12	13.34 / 14.18	15.53 / 16.17	17.78 / 18.30	19.52 / 20.13	21.37 / 22.25
	GF (D)	7.98 / <b>11.74</b>	10.91 / <b>14.49</b>	<b>14.14</b> / <b>17.09</b>	<b>17.20</b> / <b>19.60</b>	<b>20.16</b> / <b>22.22</b>	<b>22.94</b> / <b>24.84</b>	<b>25.57</b> / <b>27.6</b>
<i>Barbara</i> (512 × 512)	SF	11.27 / 14.11	-	12.71 / 14.82	-	13.97 / 15.98	-	13.39 / 16.00
	GF (S)	<b>14.56</b> / <b>14.61</b>	<b>15.71</b> / 15.10	<b>16.54</b> / 15.58	<b>17.23</b> / 16.17	18.06 / 17.10	19.06 / 18.44	20.94 / 20.95
	GF (D)	8.51 / 14.45	10.38 / <b>15.49</b>	12.50 / <b>16.80</b>	15.79 / <b>18.69</b>	<b>18.64</b> / <b>21.22</b>	<b>21.95</b> / <b>23.87</b>	<b>24.85</b> / <b>26.43</b>
<i>Boat</i> (512 × 512)	SF	14.13 / 16.16	-	16.17 / 18.04	-	17.84 / 19.91	-	17.53 / 20.13
	GF (S)	<b>16.28</b> / 17.09	<b>17.70</b> / 18.19	<b>19.21</b> / 19.53	<b>20.83</b> / 21.02	<b>22.54</b> / 22.69	24.28 / 24.35	25.99 / 26.00
	GF (D)	12.72 / <b>17.82</b>	14.85 / <b>19.44</b>	16.42 / <b>21.33</b>	19.16 / <b>23.26</b>	22.39 / <b>25.16</b>	<b>25.07</b> / <b>27.10</b>	<b>27.37</b> / <b>28.86</b>
<i>Hill</i> (512 × 512)	SF	12.89 / 16.50	-	15.10 / 17.68	-	16.39 / 18.93	-	15.63 / 18.96
	GF (S)	<b>16.28</b> / 17.34	<b>17.74</b> / 18.34	<b>18.96</b> / 19.26	<b>20.33</b> / 20.50	<b>21.62</b> / 21.72	23.27 / 23.29	24.51 / 24.51
	GF (D)	8.51 / <b>17.90</b>	10.15 / <b>19.34</b>	12.45 / <b>20.99</b>	16.33 / <b>22.89</b>	19.85 / <b>24.75</b>	<b>23.88</b> / <b>26.79</b>	<b>26.71</b> / <b>28.64</b>
<i>Lena</i> (512 × 512)	SF	13.82 / 18.78	-	16.56 / 20.52	-	18.02 / 22.27	-	18.18 / 22.56
	GF (S)	<b>18.03</b> / 19.55	<b>19.63</b> / 20.69	<b>21.32</b> / 22.11	<b>23.00</b> / 23.56	<b>24.75</b> / 25.14	<b>26.45</b> / 26.69	27.92 / 28.10
	GF (D)	11.36 / <b>20.25</b>	13.04 / <b>21.99</b>	15.49 / <b>23.95</b>	18.38 / <b>25.89</b>	21.35 / <b>27.94</b>	25.40 / <b>29.88</b>	<b>28.10</b> / <b>31.73</b>
<i>Man</i> (512 × 512)	SF	13.03 / 16.33	-	15.47 / 17.98	-	16.96 / 19.61	-	16.50 / 19.83
	GF (S)	<b>15.95</b> / 16.97	<b>17.41</b> / 18.01	<b>18.78</b> / 19.20	<b>20.27</b> / 20.48	21.76 / 21.92	23.46 / 23.57	24.97 / 25.08
	GF (D)	11.33 / <b>17.42</b>	13.97 / <b>19.03</b>	16.56 / <b>20.89</b>	18.80 / <b>22.89</b>	<b>21.94</b> / <b>25.03</b>	<b>24.83</b> / <b>27.11</b>	<b>27.65</b> / <b>29.09</b>

\* This parameter is used for GF with the constant number of acquisitions  $L = 8$ .

TABLE III  
RATE-DISTORTION PERFORMANCE OF GF WITH (D) AND WITHOUT (S) FINITE DIFFERENTIATION COMPARED TO SF [8].



Fig. 10. Reconstruction of *Bird* (256 × 256) at 1/8 bpp ( $M = 8,192$ ) using three distinct methods. From left to right: GF using  $M_0^2 = 256^2$ ,  $L = 8$ , and  $\Lambda = 1/64$  without (SNR: 21.65 dB, BSNR: 23.4 dB) and with finite differentiation (SNR: 19.79 dB, BSNR: 25.65 dB), and JPEG (SNR: 19.68 dB, BSNR: 22.66 dB). The plain-JPEG compression is performed at its lowest quality settings, which approximately yields the same bitrate (the corresponding file size is 10,280 bits, including header data).

$$\begin{aligned}
 a_2'' &= \left\{ a_2 \in \mathbb{R}_+^* : D((\cdot - \tilde{g}^{(n)})^{-2} P_2(\cdot)) = 0 \right\} \\
 &= \frac{1}{3} M \frac{(2u+1)^2}{(u^2+u+1)^2}. \quad (45)
 \end{aligned}$$

$$\begin{aligned}
 a_2 &= \max(a_2', a_2'') \\
 &= \begin{cases} M \frac{(2u+1)^2}{3(u^2+u+1)^2}, & 0 \leq u \leq 1 \\ M \frac{u(u^2+2u+3)^2}{4(u^2+u+1)^3}, & u > 1. \end{cases} \quad (46)
 \end{aligned}$$

In this first case, the three coefficients are thus determined by combining (43) and (46) given  $u$ .

Given its definition, the function  $\psi$  corresponds to the maximum between its linear and nonlinear constituents. This determines our overall first-case solution as

2) *The point  $\tilde{g}^{(n)}$  is in the linear part of  $\psi(\gamma t)$* : In this case, where  $u < 0$ , the coefficients  $a_0$  and  $a_1$  are expressed as

$$\begin{aligned} a_0 &= M^{-1}(M^{-1}u^2a_2 + 1), \\ a_1 &= -\gamma(2M^{-1}ua_2 + 1), \end{aligned} \quad (47)$$

the optimal parabola being always tangent at some distinct point in the nonlinear part of  $\psi$ . Since the corresponding polynomial  $\mathcal{P}_2$  contains one single double root in that configuration, the corresponding solution is

$$a_2 = \{a_2 \in \mathbb{R}_+^* : D(P_2(\cdot)) = 0\}. \quad (48)$$

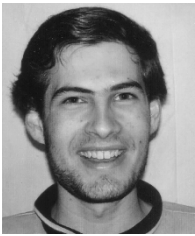
The scalar value  $a_2$  corresponds to the positive and real root of the cubic polynomial

$$\begin{aligned} \mathcal{P}_3(t) &= 12(u^2 + u + 1)^3 t^3 \\ &+ (3u^5 + 68u^4 + 214u^3 - 24u^2 - 89u + 8)Mt^2 \\ &+ (14u^3 + 168u^2 - 66u - 4)M^2t \\ &+ 27M^3u, \end{aligned} \quad (49)$$

for which the analytical expression can be found [35]. The behavior of  $\mathcal{P}_3$  as a function of  $u < 0$  guarantees the uniqueness of the solution. The coefficients are obtained in this case by solving (49) and then using (47).

#### REFERENCES

- [1] J. Romberg, "Sensing by random convolution," in *2nd IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing*, St. Thomas, United States Virgin Islands, December 12-14, 2007, pp. 137–140.
- [2] A. Stern, Y. Rivenson, and B. Javidi, "Optically compressed image sensing using random aperture coding," in *Proceedings of the International Society for Optical Engineering*, Orlando, FL, USA, March 17-18, 2008, vol. 6975, pp. 69750D–1–10.
- [3] M.F. Duarte, M.A. Davenport, D. Takhar, J. Laska, T. Sun, K.F. Kelly, and R.G. Baraniuk, "Single-pixel imaging via compressive sampling: Building simpler, smaller, and less-expensive digital cameras," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 83–91, March 2008.
- [4] R.F. Marcia and R.M. Willett, "Compressive coded aperture super-resolution image reconstruction," in *IEEE International Conference on Acoustic, Speech and Signal Processing*, Las Vegas, NV, USA, March 31-April 4, 2008, pp. 833–836.
- [5] F. Seibert, Y.M. Zou, and L. Ying, "Toeplitz block matrices in compressed sensing and their applications in imaging," in *Proceedings of the 5th International Conference on Information Technology and Application in Biomedicine*, Shenzhen, China, May 30-31, 2008, pp. 47–50.
- [6] A.M. Bruckstein, D.L. Donoho, and M. Elad, "From sparse solutions of systems of equations to sparse modeling of signals and images," *SIAM Review*, vol. 51, no. 1, pp. 34–81, February 2009.
- [7] P.T. Boufounos and R.G. Baraniuk, "1-Bit compressive sensing," in *42nd Annual Conference on Information Sciences and Systems*, Princeton, NJ, USA, March 19-21, 2008, pp. 16–21.
- [8] A. Bourquard, F. Aguet, and M. Unser, "Optical imaging using binary sensors," *Optics Express*, vol. 18, no. 5, pp. 4876–4888, March 2010.
- [9] L. Jacques, J.N. Laska, P.T. Boufounos, and R.G. Baraniuk, "Robust 1-bit compressive sensing via binary stable embedding of sparse vectors," arXiv:submit/0417664 [cs.IT], February 2012.
- [10] Y. Plan and R. Vershynin, "One-Bit compressed sensing by linear programming," arXiv:1109.4299v4 [cs.IT], October 2011.
- [11] M. Bigas, E. Cabruja, J. Forest, and J. Salvi, "Review of CMOS image sensors," *Microelectronics Journal*, vol. 37, no. 5, pp. 433–451, May 2006.
- [12] J.W. Goodman, *Introduction to Fourier Optics*, McGraw Hill Higher Education, 2nd edition, 1996.
- [13] M. Unser, "Splines: A perfect fit for signal and image processing," *IEEE Signal Processing Magazine*, vol. 16, no. 6, pp. 22–38, November 1999.
- [14] W.U. Bajwa, J.D. Haupt, G.M. Raz, S.J. Wright, and R.D. Nowak, "Toeplitz-structured compressed sensing matrices," in *IEEE Workshop on Statistical Signal Processing Proceedings*, Madison, WI, United States, August 26-29, 2007, pp. 294–298.
- [15] H. Rauhut, "Circulant and Toeplitz matrices in compressed sensing," in *Proceedings of SPARS'09*, Saint-Malo, France, April 6-9, 2009.
- [16] J. Romberg and R. Neelamani, "Sparse channel separation using random probes," *Inverse Problems*, vol. 26, no. 11, pp. 115015 (25 pp.), November 2010.
- [17] J.P. Slavinsky, J.N. Laska, M.A. Davenport, and R.G. Baraniuk, "The compressive multiplexer for multi-channel compressive sensing," in *Proceedings of ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing*, Prague, Czech Republic, May 22, 2011, number 5947224 in 15206149, pp. 3980–3983.
- [18] J.N. Laska, Z. Wen, W. Yin, and R.G. Baraniuk, "Trust, but verify: Fast and accurate signal recovery from 1-bit compressive measurements," *IEEE Transactions on Signal Processing*, vol. 59, no. 11, pp. 5289–5301, November 2011.
- [19] Y. Plan and R. Vershynin, "Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach," arXiv:1202.1212v1 [cs.IT], February 2012.
- [20] Y.E. Nesterov, "A method of solving a convex programming problem with convergence speed  $O(1/k^2)$ ," *Doklady Akademii Nauk SSSR*, vol. 27, no. 2, pp. 372–376, 1983.
- [21] R.T. Rockafellar, *Convex Analysis (Princeton Mathematical Series)*, Princeton University Press, 1970.
- [22] M. Elad, P. Milanfar, and R. Rubinfeld, "Analysis versus synthesis in signal priors," *Inverse Problems*, vol. 23, no. 3, pp. 947–968, June 2007.
- [23] I.W. Selesnick and M.A.T. Figueiredo, "Signal restoration with overcomplete wavelet transforms: Comparison of analysis and synthesis priors," in *Proceedings of the SPIE - The International Society for Optical Engineering*, San Diego, CA, USA, August 2–4, 2009, p. 74460D (15 pp.).
- [24] L. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D*, vol. 60, no. 1-4, pp. 259–268, November 1992.
- [25] E.J. Candès and T. Tao, "Near-Optimal signal recovery from random projections: Universal encoding strategies?," *IEEE Transactions on Information Theory*, vol. 52, no. 12, pp. 5406–5425, December 2006.
- [26] M.M. Marim, E.D. Angelini, and J.-C. Olivo-Marin, "A compressed sensing approach for biological microscopy image denoising," in *2009 IEEE International Symposium on Biomedical Imaging: From Nano to Macro (ISBI)*, Boston, MA, USA, June 18–July 1, 2009, pp. 1374–1377.
- [27] S. Becker, J. Bobin, and E.J. Candès, "NESTA: A fast and accurate first-order method for sparse recovery," *SIAM Journal on Imaging Sciences*, vol. 4, no. 1, pp. 1–39, 2011.
- [28] M. Nikolova and K. Michael, "Analysis of half-quadratic minimization methods for signal and image recovery," *SIAM Journal of Scientific Computing*, vol. 27, no. 3, pp. 937–966, December 2005.
- [29] S. McKee, M.F. Tomé, V.G. Ferreira, J.A. Cuminato, A. Castelo, F.S. Sousa, and N. Mangiavacchi, "The MAC method," *Computers and Fluids*, vol. 37, no. 8, pp. 907–930, September 2008.
- [30] J.P. Oliveira, J.M. Bioucas-Dias, and M.A.T. Figueiredo, "Adaptive total variation image deblurring: A majorization-minimization approach," *Signal Processing*, vol. 89, no. 9, pp. 1683–1693, September 2009.
- [31] R. Pan and S.J. Reeves, "Efficient Huber-Markov edge-preserving image restoration," *IEEE Transactions on Image Processing*, vol. 15, no. 12, pp. 3728–3735, December 2006.
- [32] R. Chan, J.G. Nagy, and R.J. Plemmons, "FFT-Based preconditioners for Toeplitz-block least squares problems," *SIAM Journal on Numerical Analysis*, vol. 30, no. 6, pp. 1740–1768, December 1993.
- [33] G.E.P. Box and G. Jenkins, *Time Series Analysis: Forecasting and Control*, Holden-Day, 1976.
- [34] A. Schulz, L. Velho, and E.A.B. da Silva, "On the empirical rate-distortion performance of compressive sensing," in *Proceedings of the 2009 16th IEEE International Conference on Image Processing*, Cairo, Egypt, November 7–12, 2009, pp. 3049–3052.
- [35] G. Cardano, *Artis Magnae, Sive de Regulis Algebraicis Liber Unus*, Nuremberg, 1545.



**Aurélien Bourquard** was born in Switzerland on February 7, 1985. He received his M.Sc. degree in Microengineering from the Swiss Federal Institute of Technology (EPFL), Lausanne, Switzerland. Currently, he is a Ph.D. student with the Biomedical Imaging Group (BIG), EPFL. His research interests include image reconstruction using convex optimization and multigrid techniques, as well as new acquisition methods in the framework of computational optics.



**Michael Unser** (M'89-SM'94-F'99) received the M.S. (summa cum laude) and Ph.D. degrees in Electrical Engineering in 1981 and 1984, respectively, from the École polytechnique fédérale de Lausanne (EPFL), Switzerland. From 1985 to 1997, he worked as a scientist with the National Institutes of Health, Bethesda USA. He is now full professor and Director of the Biomedical Imaging Group at the EPFL.

His main research area is biomedical image processing. He has a strong interest in sampling theories, multiresolution algorithms, wavelets, and the use of splines for image processing. He has published 200 journal papers on those topics, and is one of ISI's Highly Cited authors in Engineering (<http://isihighlycited.com>).

Dr. Unser has held the position of associate Editor-in-Chief (2003-2005) for the IEEE Transactions on Medical Imaging and has served as Associate Editor for the same journal (1999-2002; 2006-2007), the IEEE Transactions on Image Processing (1992-1995), and the IEEE Signal Processing Letters (1994-1998). He is currently member of the editorial boards of Foundations and Trends in Signal Processing, and Sampling Theory in Signal and Image Processing. He co-organized the first IEEE International Symposium on Biomedical Imaging (ISBI2002) and was the founding chair of the technical committee of the IEEE-SP Society on Bio Imaging and Signal Processing (BISP).

Dr. Unser received the 1995 and 2003 Best Paper Awards, the 2000 Magazine Award, and the 2008 Technical Achievement Award from the IEEE Signal Processing Society. He is an EURASIP Fellow and a member of the Swiss Academy of Engineering Sciences.